

1) A) (i) Label : 1

(ii)

		Predicted	
		+	-
Actual	+	TP	FN
	-	FP	TN

		Predicted	
C	O	+	-
O	+	209	64
V	-	134	93
G	-		

(iii) TP Rati =  $\frac{\# TP}{\# TP + \# FN} = \frac{209}{209 + 64}$

$$= \frac{209}{273} = 0.766$$

FP Rati =  $\frac{\# FP}{\# FP + \# TN} = \frac{134}{134 + 93}$

$$= \frac{134}{227} = 0.590$$

$$\text{Accuracy} = \frac{\text{Correct Predicted}}{\text{Total Predicted}} = \frac{209 + 93}{500} \times 100$$

$$= (0.604) \times 100 = \underline{\underline{60.4\%}}$$

(iv) Label: 1

Q7

		Predicted	
		+	-
S	+	212	61
	-	136	91

$$(vi) \quad \text{TP Ratio} = \frac{\# \text{TP}}{\# \text{TP} + \# \text{FN}} = \frac{212}{212 + 61}$$

$$= \frac{212}{273} = 0.778$$

$$\text{FP Ratio} = \frac{\# \text{FP}}{\# \text{FP} + \# \text{TN}} = \frac{136}{136 + 91} = \frac{136}{227}$$

$$= 0.599$$

$$\text{Accuracy} = 100 \times \frac{\text{Correct Predicted}}{\text{Total Predicted}} = 100 \times \frac{212 + 91}{500}$$

$$= 100 \times \frac{303}{500} = 0.606 \times 100 = 60.6\%.$$

Viii)

		Predicted	
		+	-
Actual	+	273	6
	-	227	0

(b)

Distance function that we are using would not be a good choice. Because in this we are just looking at the unique tokens between 2 comment as a measure for distance & then dividing it by 1.

But this would give ~~rise~~ <sup>certain</sup> to certain problems

1) Case 1:- When  $C_1$  is small as compared to  $C_2$ , let  $C_1$  &  $C_2$  both belongs to the class.

Even if they both belong to same class they should be at a long distance from each other. Because they would only have few tokens in common. As compared to other example in which  $C_1$  &  $C_2$  are of same length. So  $C_2$  would not be predicted as +ve.

Case 2: When  $C_1$  &  $C_2$  are both of large length. Let  $C_1$  belongs to -ve Class &  $C_2$  to +ve.

Then ~~the~~ distance between them would be less as they would have a lot of ~~tokens~~ common tokens (like Punctuation, 'am', 'is', etc)

in Common. So  $C_2$  they would be predicted as ' $\sqrt{a}$ ' even when it is 'two'.

- 2) We are dividing 1 by the unique taken in  $C_1$  &  $C_2$ .  
So when they both have large Number of  
factors in Common. So dividing that by  
1 would make it pretty small. So a lot  
of examples would be stored ~~not~~ calculated  
to be <sup>near</sup> <sub>near</sub> each other as the difference would  
be in some <sup>small</sup> decimal Number between  
them. If the difference is only at certain  
decimal point like  $0.0017893631$  &  $0.0017893630$   
then it would be hard for us to predict  
accurately because we store number as an  
approximation in Variables.

(d) Accuracy for  $K=3$  is 65.9%.

Accuracy for  $K=7$  is 68.8%.

Accuracy for  $K=9$  is 61.2%.

(ii) Accuracy on test data set with  $K=3$  is 59%.

We have used  $K=3$  because it has maximum accuracy on training.

		Predicted	
		+	-
Actual	+	212	61
	-	144	83

d (i) Distance function =  $\frac{1}{\text{uncommon unique words in } C_1 \cap C_2}$

Uncommon unique words means only rare or less common words that if we remove the words like "the", "he", "she", "are", "is", "at" etc. as well as punctuation marks like ., :, !, ; etc.

For e.g.  $C_1 =$  "Hello, its me. You look nice."  
 $C_2 =$  "nice of you., The weather."

Hence uncommon = {Nice}  
which is Pn  $C_1$  and  $C_2$

It would remove tokens like "you", ".," , Because they are common words, and hence don't influence the positivity or negativity of a comment.

(ii) In first, we calculate the number of unique words that are common in both  $C_1$  and  $C_2$ . But, when we do that we count tokens like '.', '!', 'and' etc. even when they have no influence on the positivity or negativity of a comment.

So when we remove these common tokens, we are left with exactly the number of tokens which are influential. Thus, the tokens will have a greater influence on the comment being positive or negative.

### (iii) Confusion Matrix

		Predicted	
Correct	+	TP 230	FN 43
	-	FP 88	TN 139

$$(iv) \text{ True Positive Rate} = \frac{\# TP}{\# TP + \# FN}$$

$$= \frac{230}{230 + 43} = \frac{230}{273}$$

$$\approx 0.84$$

$$\text{False Positive Rate} = \frac{\# FP}{\# FP + \# TN}$$

$$= \frac{88}{88 + 139}$$

$$= 0.388$$

$$\text{Accuracy} = \frac{\# TP + \# TN}{\text{Total}}$$

$$= \frac{230 + 139}{500}$$

$$= 0.738 \text{ or } 73.8\%$$

(v) Confusion Matrix for K=5

		Predicted	
		+	-
Actual	+	235	38
	-	94	133

vii) True Positive rate =  $\frac{\# TP}{\# TP + \# FN} = \frac{238}{235 + 38}$

$$= \frac{238}{273} = 0.872$$

False Positive rate =  $\frac{\# FP}{\# FP + \# TN} = \frac{94}{94 + 133}$

$$= \frac{94}{227} = 0.414$$

Accuracy =  $\left( \frac{TP + TN}{Total} \right) \times 100 = \left( \frac{235 + 133}{500} \right) \times 100 = 0.736 \times 100$

$$= 73.6\%$$

(vii) Yes, it did achieves higher accuracy  
both for  $k=1$  and  $k=5$

for  $k=1$

Old function accuracy = 60.4%

New distance function accuracy = 73.8%

For  $k = 5$

Old distance function accuracy = 60.6 %  
New distance function accuracy = 73.6 %

~~Anc~~